# Pandia: comprehensive contention-sensitive thread placement

**Daniel Goodman**
**Georgios Varisteas**
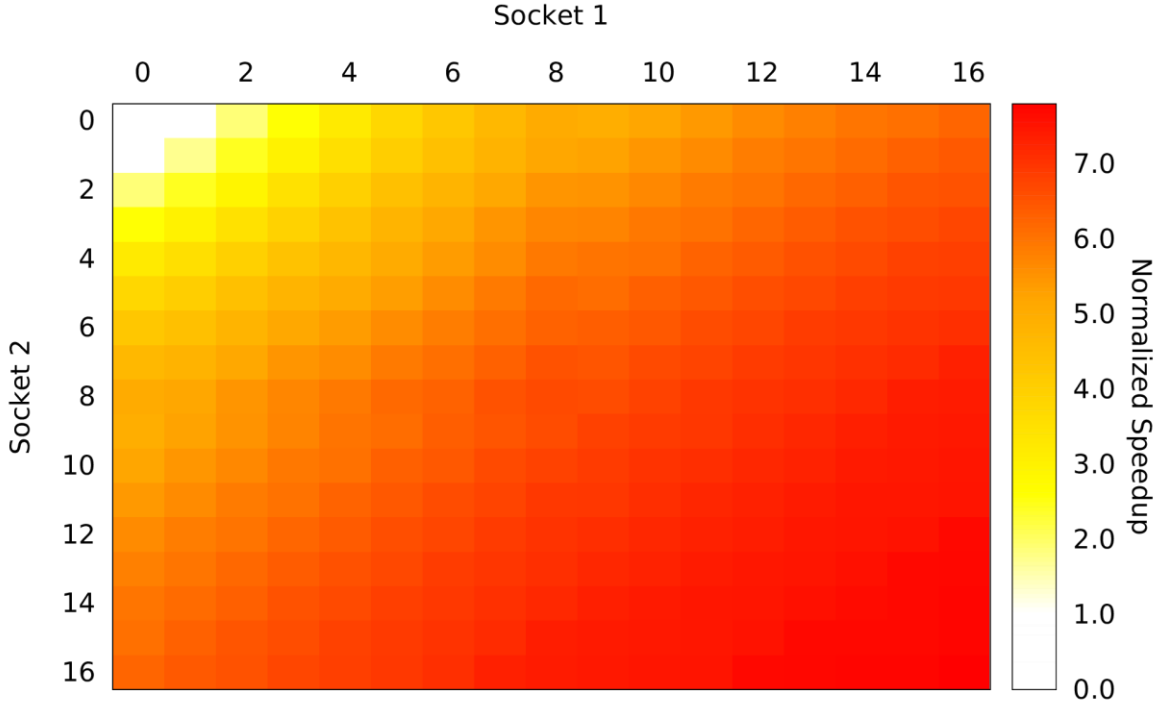**Tim Harris**

daniel.goodman@oracle.com

ORACLE®

# Safe Harbor Statement

The following is intended to outline our general product direction. It is intended for information purposes only, and may not be incorporated into any contract. It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions. The development, release, and timing of any features or functionality described for Oracle's products remains at the sole discretion of Oracle.
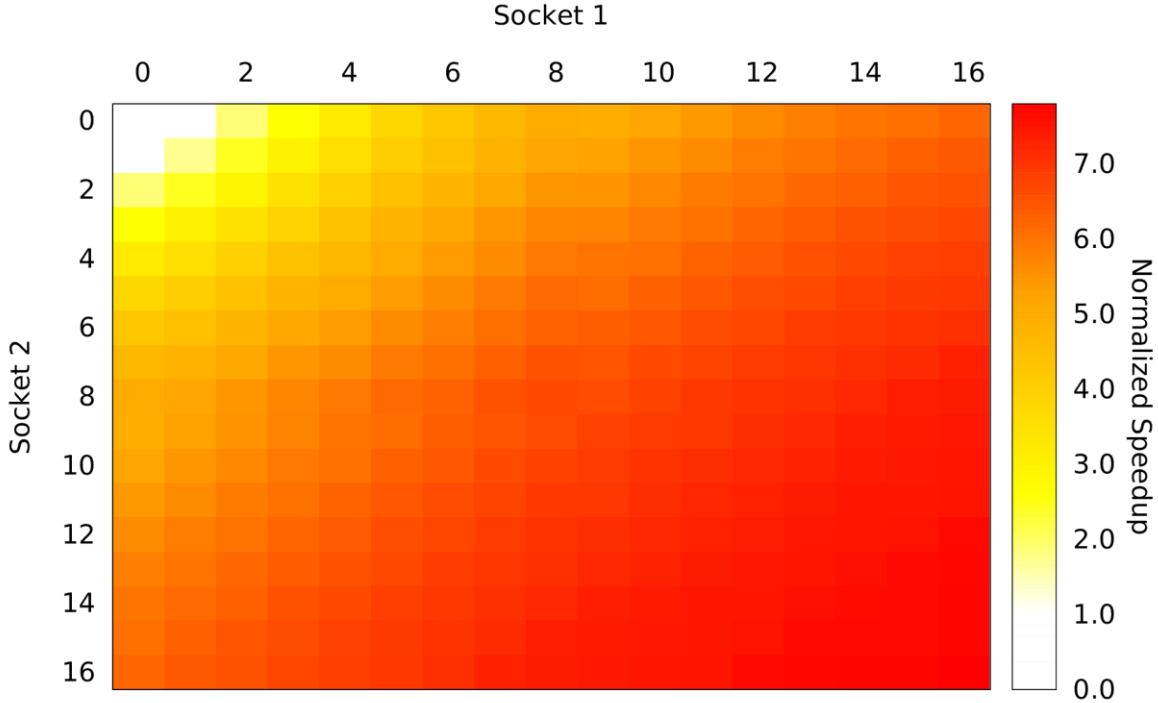
# Example behaviour patterns
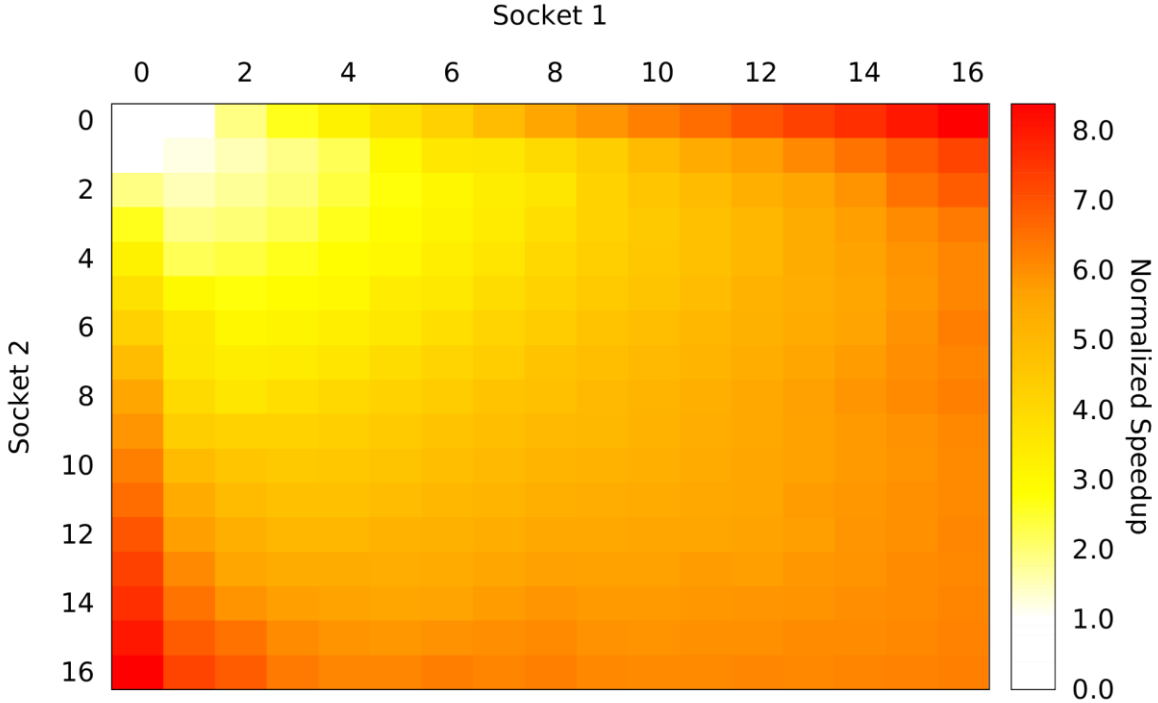
## Parallel radix join optimized

# Example behaviour patterns

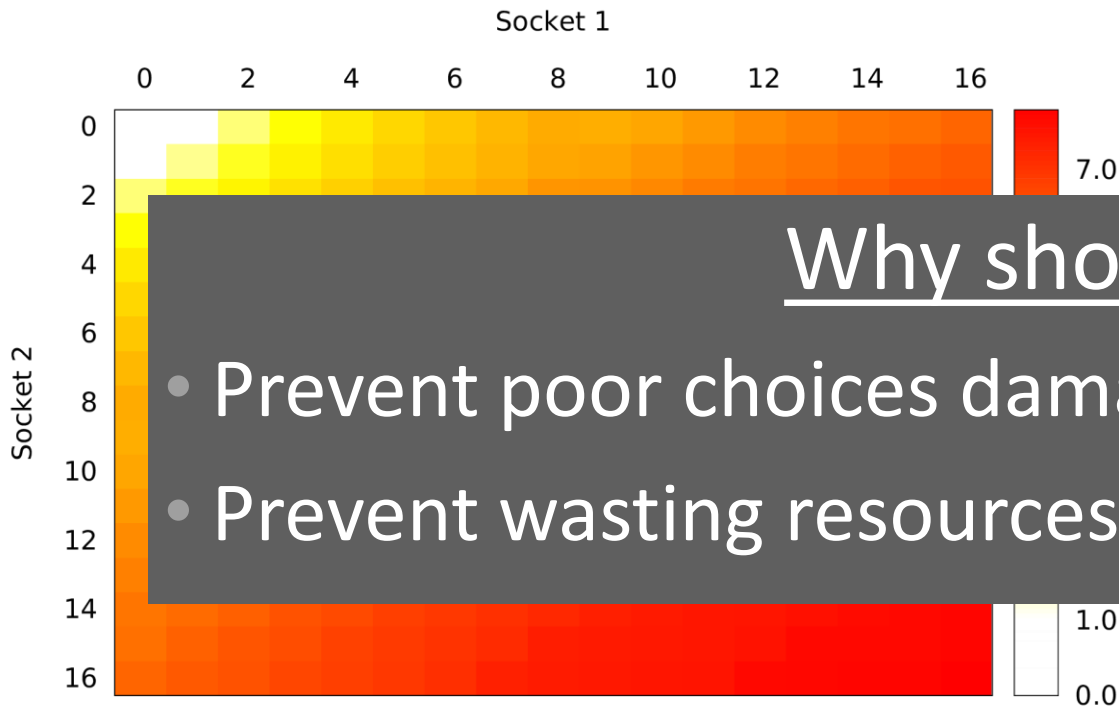## Parallel radix join optimized
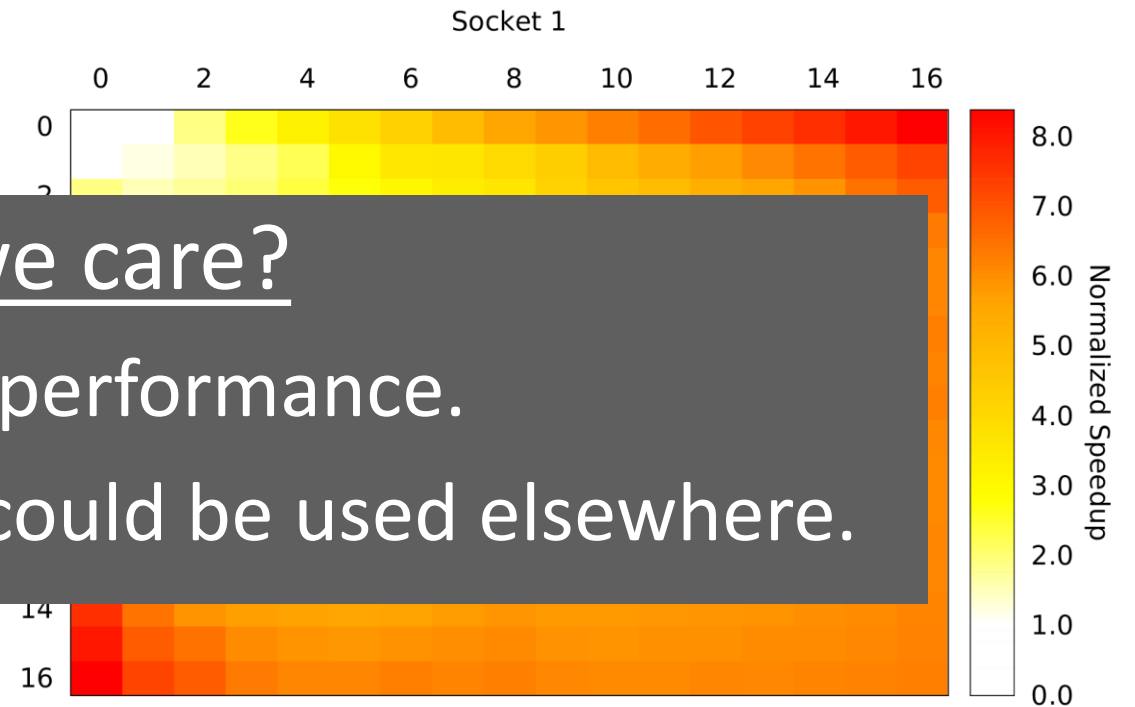


## Molecular dynamics simulation

# Example behaviour patterns

## Parallel radix join optimized

## Molecular dynamics simulation



### Why should we care?

- Prevent poor choices damaging performance.
- Prevent wasting resources that could be used elsewhere.
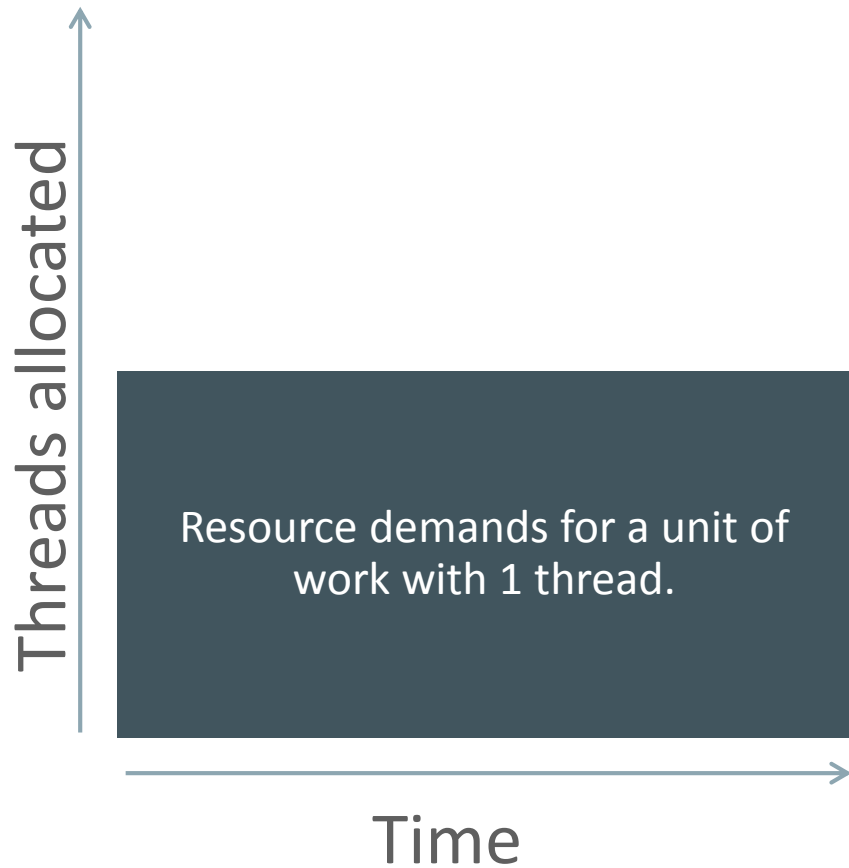
# Overview

1 What is the problem?

**2 How does Pandia work?**

3 How well does Pandia work?

4 Conclusions

ORACLE®

# Key idea 1

**Our workloads perform a roughly constant amount of work**
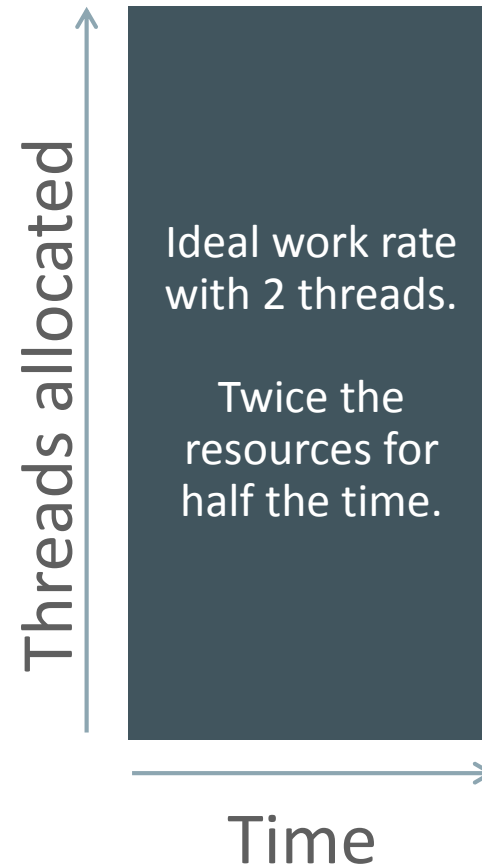
# Key idea 1

**Our workloads perform a roughly constant amount of work**



Threads allocated

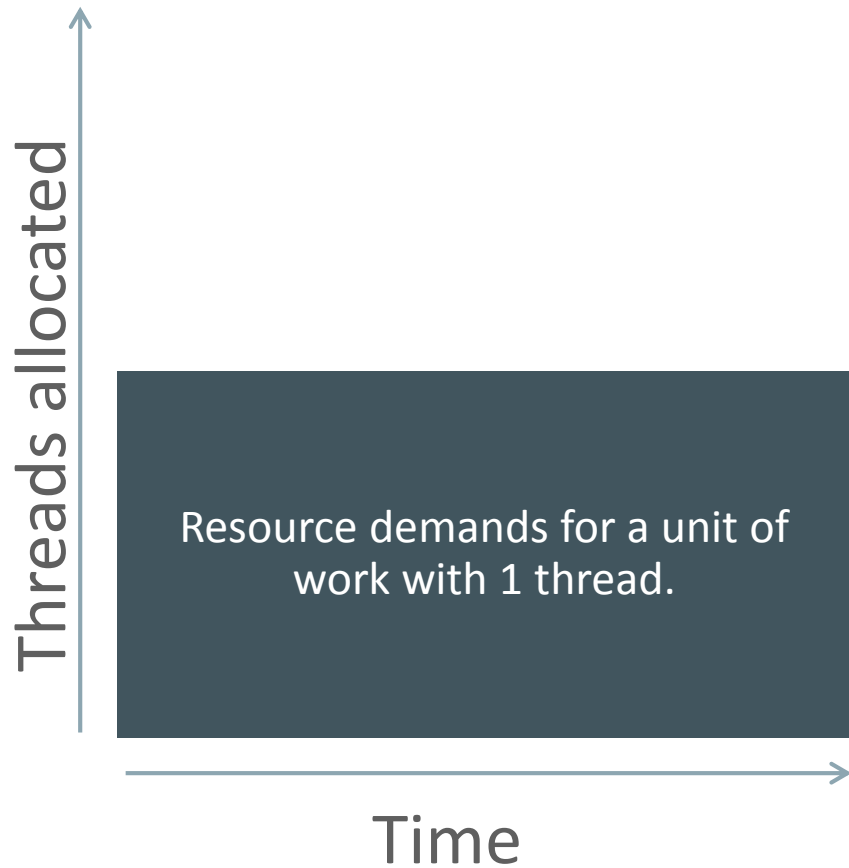Resource demands for a unit of work with 1 thread.

Time

# Key idea 1

**Our workloads perform a roughly constant amount of work**

# Key idea 1

## Our workloads perform a roughly constant amount of work
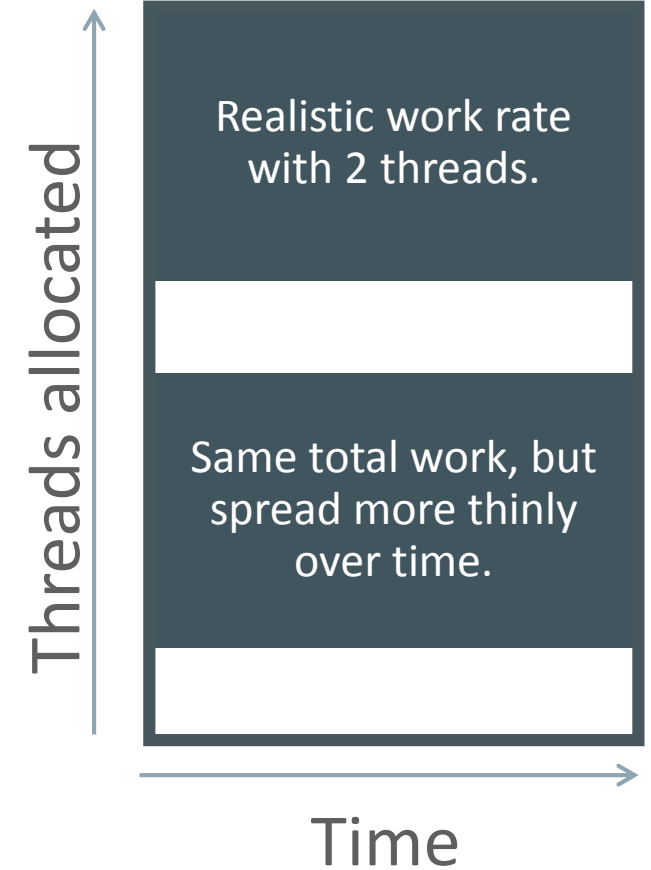


(Left chart) Threads allocated vs. Time: **Resource demands for a unit of work with 1 thread.**

(Middle chart) Threads allocated vs. Time: **Ideal work rate with 2 threads. Twice the resources for half the time.**

(Right chart) Threads allocated vs. Time: **Realistic work rate with 2 threads. Same total work, but spread more thinly over time.**

# Key idea 2

**Hardware is getting ever more complicated, but looks simpler**

- This makes it harder to model using conventional techniques, but removes much of the need to model behaviour in detail.

**ORACLE**®

# Key idea 2

**Hardware is getting ever more complicated, but looks simpler**

- This makes it harder to model using conventional techniques, but removes much of the need to model behaviour in detail.
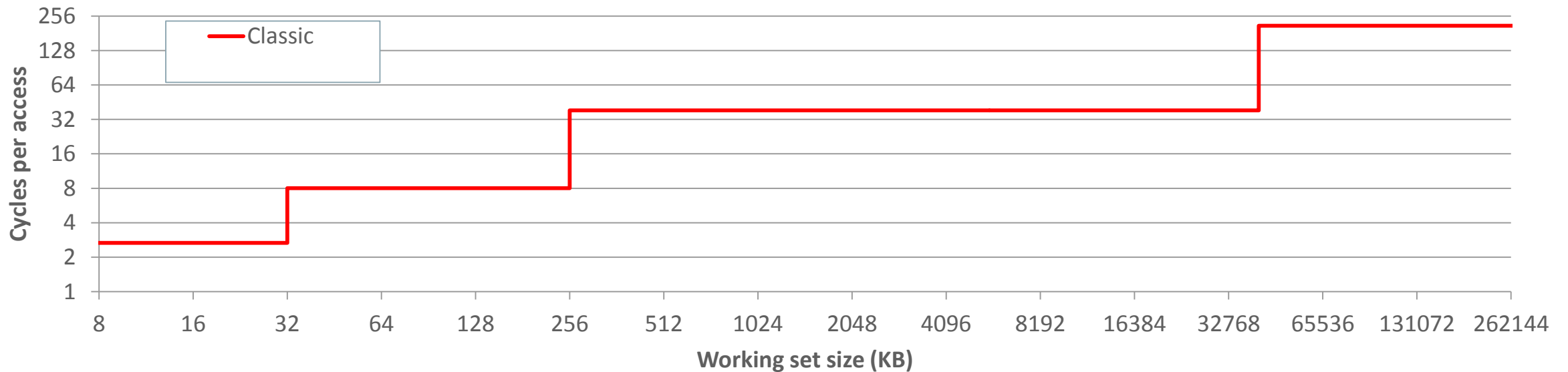
- Example: Adaptive caches.

# Key idea 2

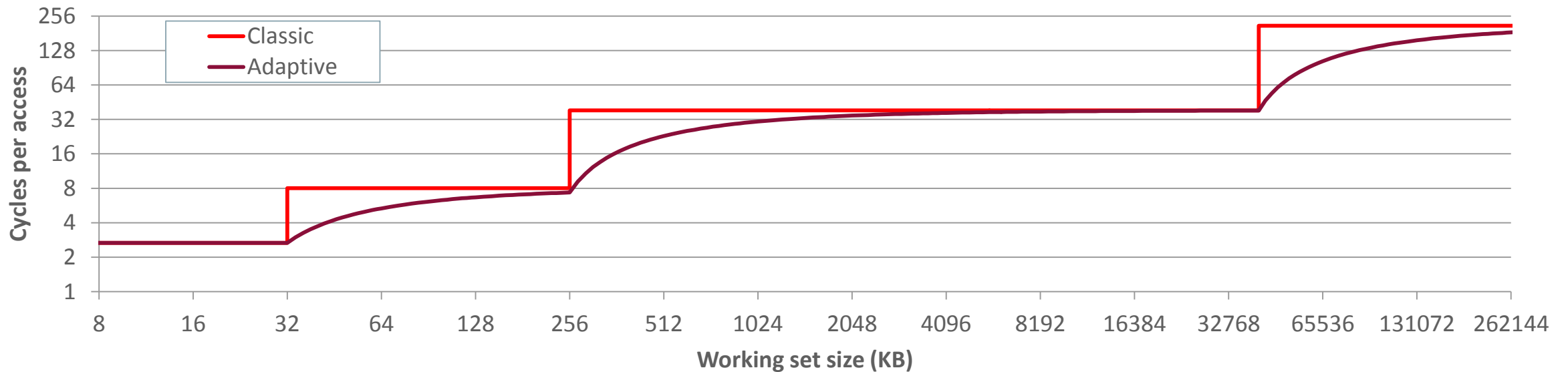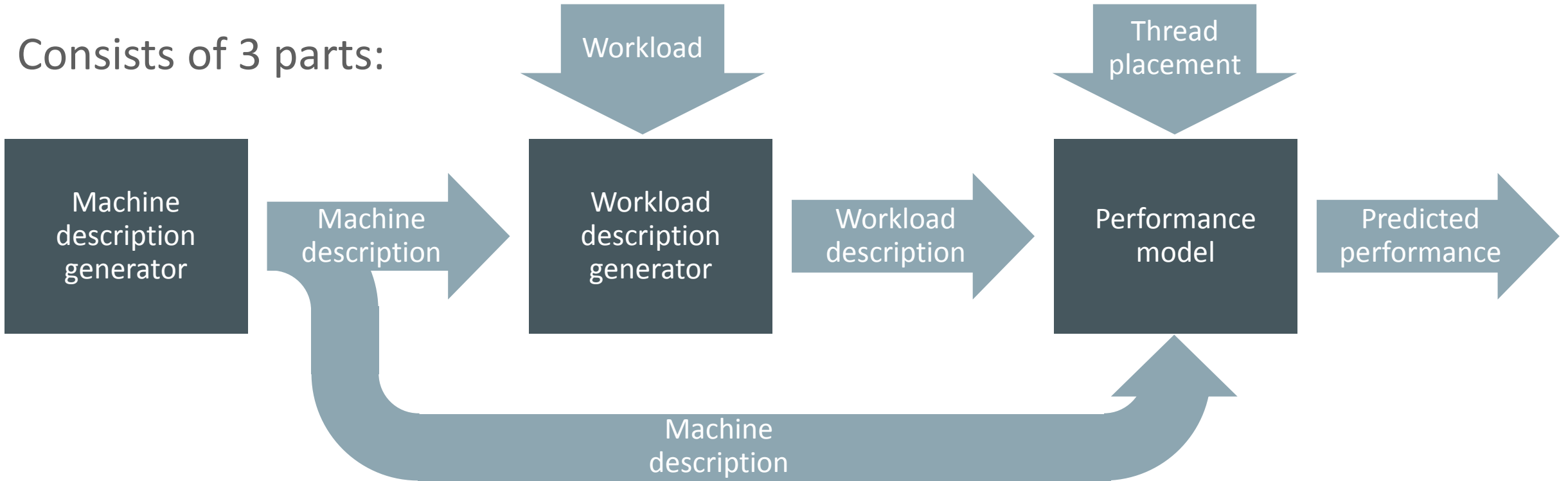**Hardware is getting ever more complicated, but looks simpler**

- This makes it harder to model using conventional techniques, but removes much of the need to model behaviour in detail.

- Example: Adaptive caches.

# Pandia: **predicting the performance of in-memory workloads**

Consists of 3 parts:

# Machine description

Machine description generator
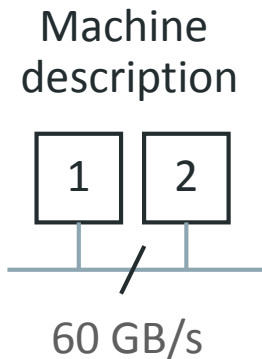
ORACLE®

# Machine description

- Query the OS/processor for CPU count, core count, cache sizes, …

Machine description



Machine description generator

# Machine description

Machine
description

**Machine
description
generator** ➡

1  2

60 GB/s

- Query the OS/processor for CPU count, core count, cache sizes, …

- Measure synthetic stress applications for:
  - Latency
  - Bandwidth
  - Execution rate (normalized IPC)
  - …

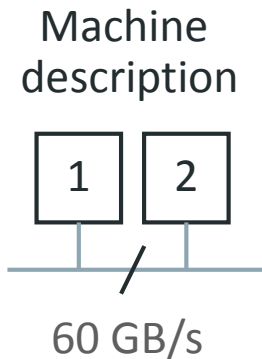**ORACLE**®

# Machine description

Machine
description

Machine
description
generator ➡

1 2

60 GB/s

- Query the OS/processor for CPU count, core count, cache sizes, …

- Measure synthetic stress applications for:
  - Latency
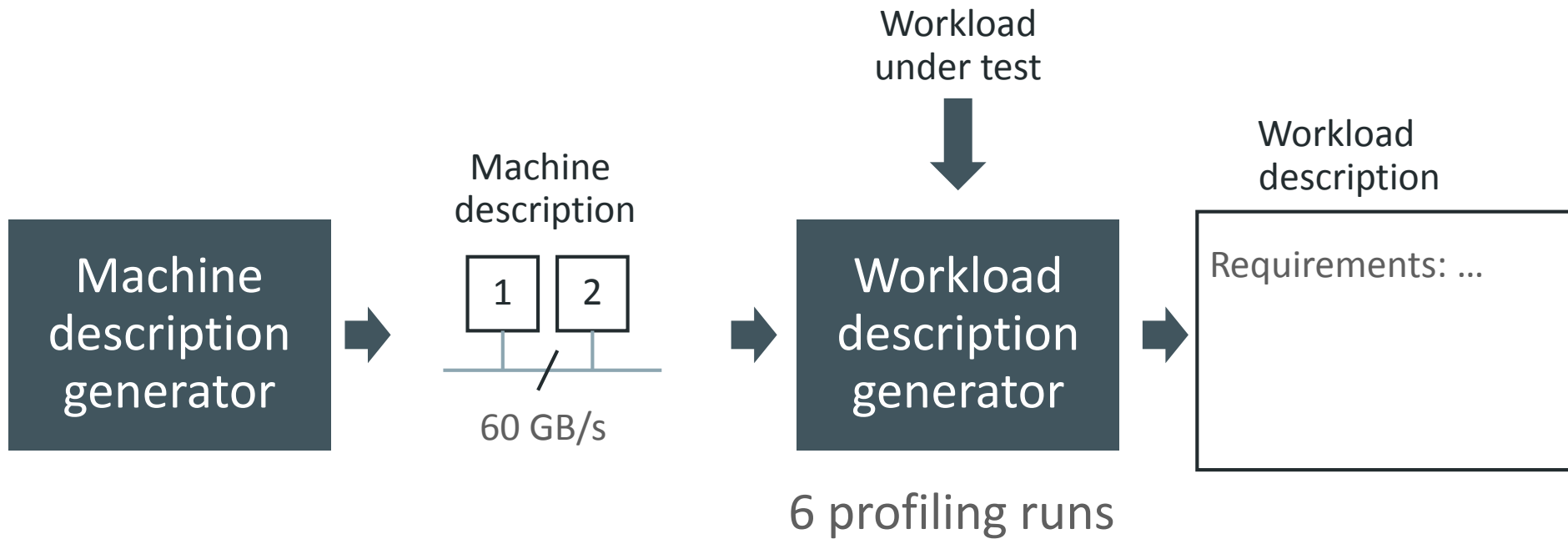  - Bandwidth
  - Execution rate (normalized IPC)
  - …

- Detailed statistics are gained from performance counters
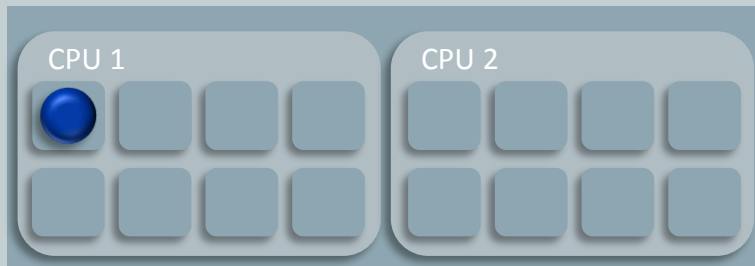
# Workload description

# Workload description: per-thread requirements

- Characteristics that reflect the requirements of an individual thread:
  - Instruction execution rate (normalized instructions per cycle – IPC)
  - Memory bandwidth
  - Inter cache bandwidth
  - …
- Measured while running the application
  - Run the application with a minimal thread count
  - Record statistics using performance counters

# Workload description: runs



Run 1: Single Thread

CPU 1

CPU 2

# Workload description: parallelism characteristics

- Characteristics of the interactions between the threads
- These all reflect synchronization either at the hardware level or within the application

ORACLE

# Workload description: parallelism characteristics

- Characteristics of the interactions between the threads
- These all reflect synchronization either at the hardware level or within the application.

| Parallelism<br>The percentage of the executed code that is parallel | |
| --- | --- |
| | |

**ORACLE**

# Workload description: runs

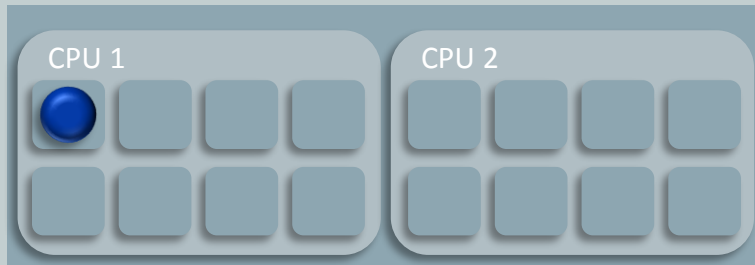| Run 1: single thread | Run 2: n threads, single socket |
|---|---|

# Workload description: parallelism characteristics

- Characteristics of the interactions between the threads
- These all reflect synchronization either at the hardware level or within the application

| Parallelism | Communication slowdown |
|---|---|
| The percentage of the executed code that is parallel | Slowdown due to latency of thread communications between sockets, nodes, … |
| | |

# Workload description: runs

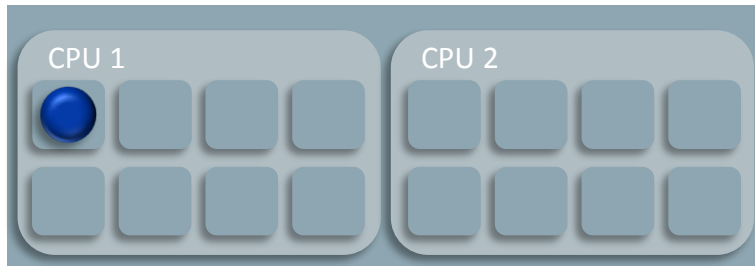| Run 1: single thread | Run 2: n threads, single socket | Run 3: symmetric, multi-socket |
|---|---|---|

# Workload description: parallelism characteristics

- Characteristics of the interactions between the threads
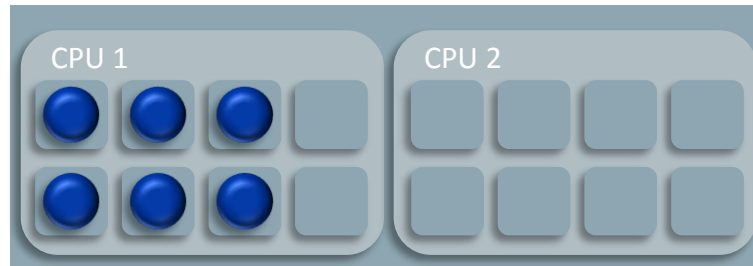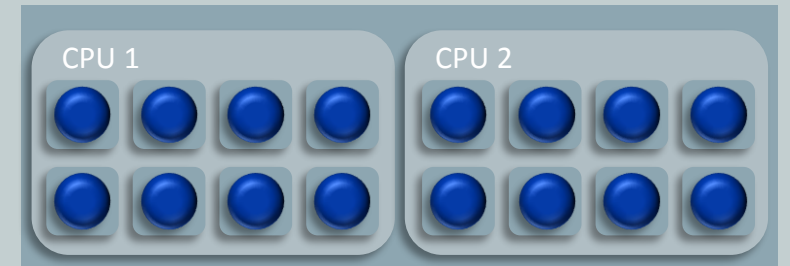- These all reflect synchronization either at the hardware level or within the application

| | |
|---|---|
| **Parallelism**<br>The percentage of the executed code that is parallel | **Communication slowdown**<br>Slowdown due to latency of thread communications between sockets, nodes, … |
| **Thread interlocking**<br>Are threads independent? | |

ORACLE®

# Workload description: runs

**Run 1: single thread**



**Run 2: n threads, single socket**



**Run 3: symmetric, multi-socket**



**Run 4: n threads, n slowed**
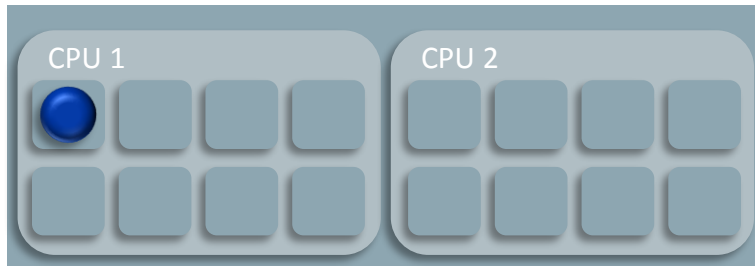


**Run 5: n threads, 1 slowed**

# Workload description: parallelism characteristics

- Characteristics of the interactions between the threads
- These all reflect synchronization either at the hardware level or within the application

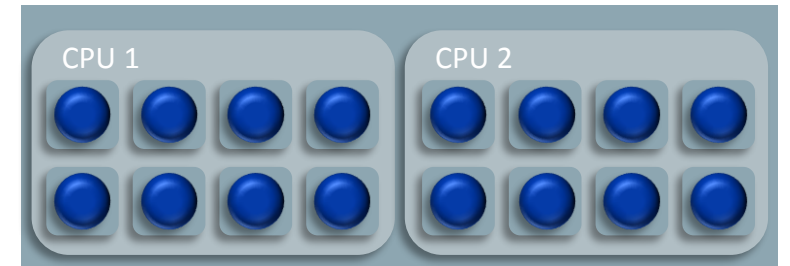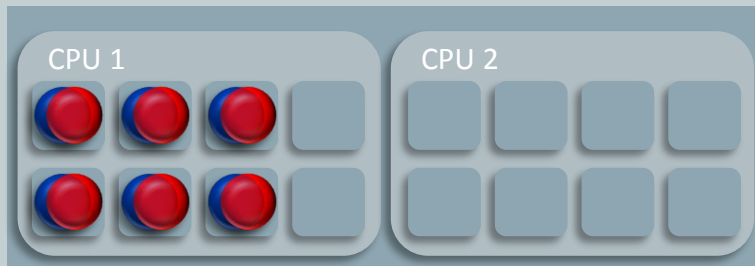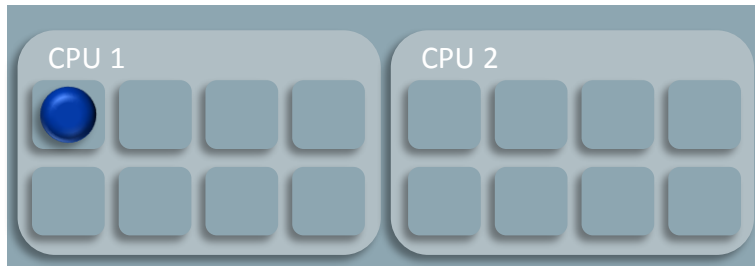| **Parallelism** <br> The percentage of the executed code that is parallel | **Communication slowdown** <br> Slowdown due to latency of thread communications between sockets, nodes, … |
|---|---|
| **Thread interlocking** <br> Are threads independent? | **Coincident resource demands** <br> Do threads have synchronized mode changes (e.g., intensive reads followed by intensive CPU) |

ORACLE®

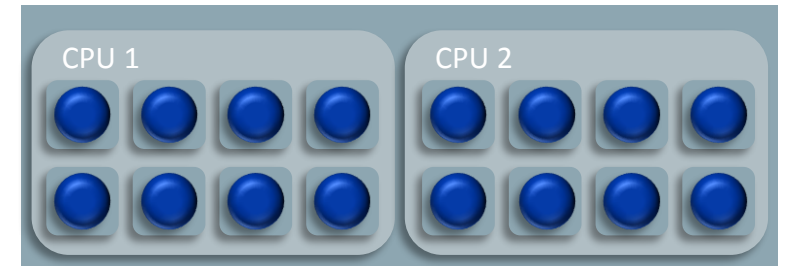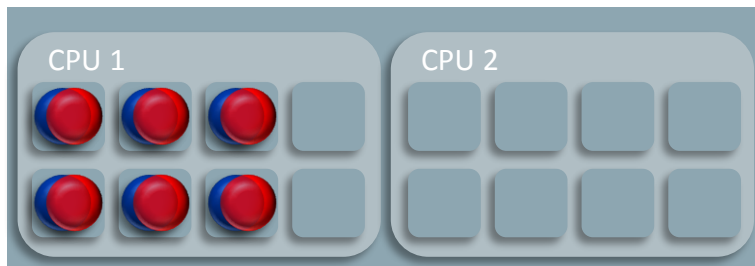Workload description: runs

Run 1: Single thread

Run 2: n threads, single socket

Run 3: symmetric, multi-socket

Run 4: n threads, n slowed

Run 5: n threads, 1 slowed

Run 6: n threads, hyper-threaded

# Performance prediction

# Performance model

**Predicted speedup** $=$ $\dfrac{\text{Estimated speedup with thread count (Amdahl's law)}}{\text{Estimated slowdown due to resource contention}}$

# Estimating thread slowdown

- A machine is a set of components each with a set of resources

- Over subscribing any of these resources will produce a slowdown

- The nature of the slowdown will vary with the resource

Max

0

NORMALIZED IPC

ORACLE®

# Estimating thread slowdown

- A machine is a set of components each with a set of resources

- Over subscribing any of these resources will produce a slowdown

- The nature of the slowdown will vary with the resource

Max

NORMALIZED IPC

0

ORACLE®

# Estimating thread slowdown

- A machine is a set of components each with a set of resources

- Over subscribing any of these resources will produce a slowdown

- The nature of the slowdown will vary with the resource
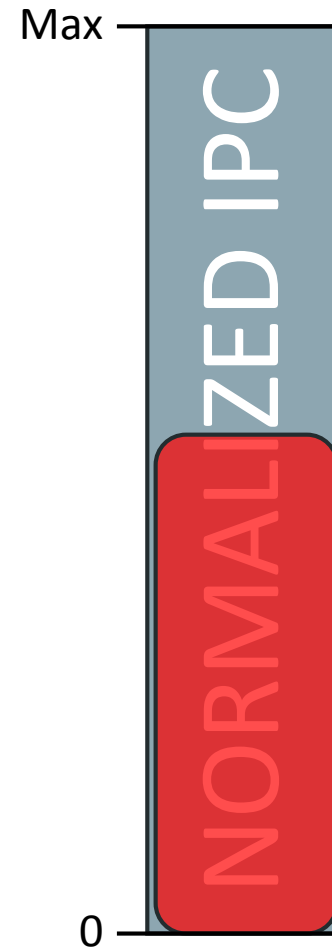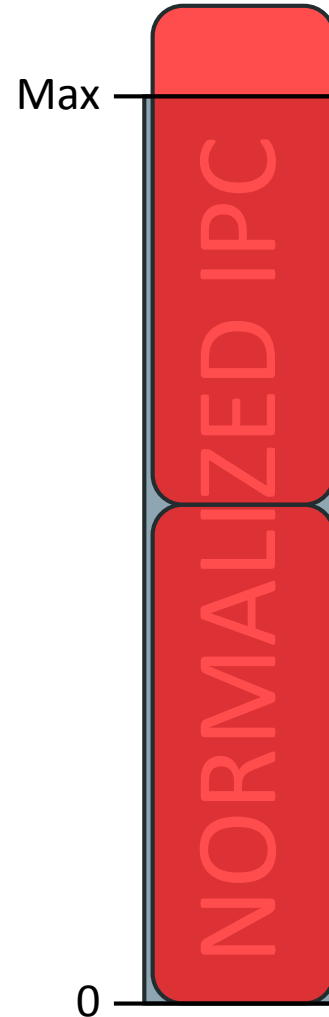
Max

NORMALIZED IPC

0

ORACLE®

# Estimating thread slowdown

- A machine is a set of components each with a set of resources

- Over subscribing any of these resources will produce a slowdown

- The nature of the slowdown will vary with the resource
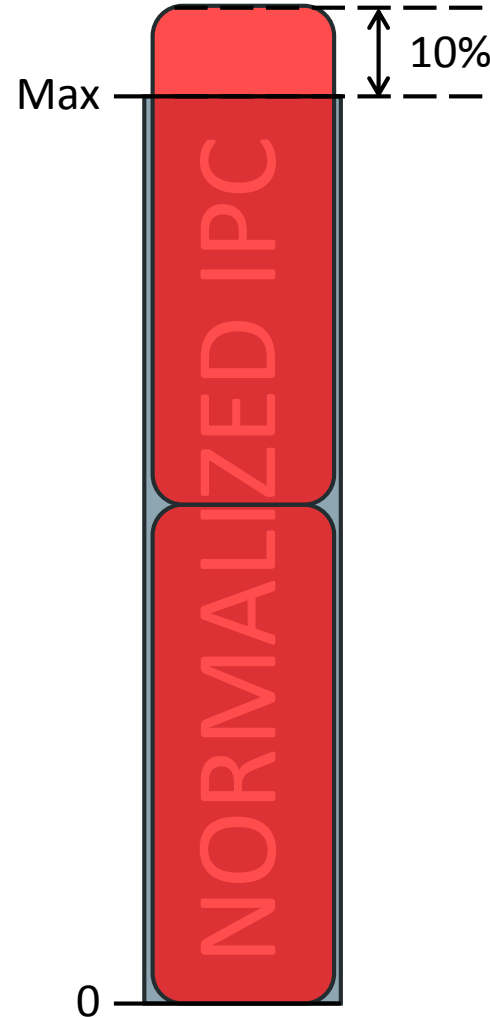


10%

Max

0

NORMALIZED IPC
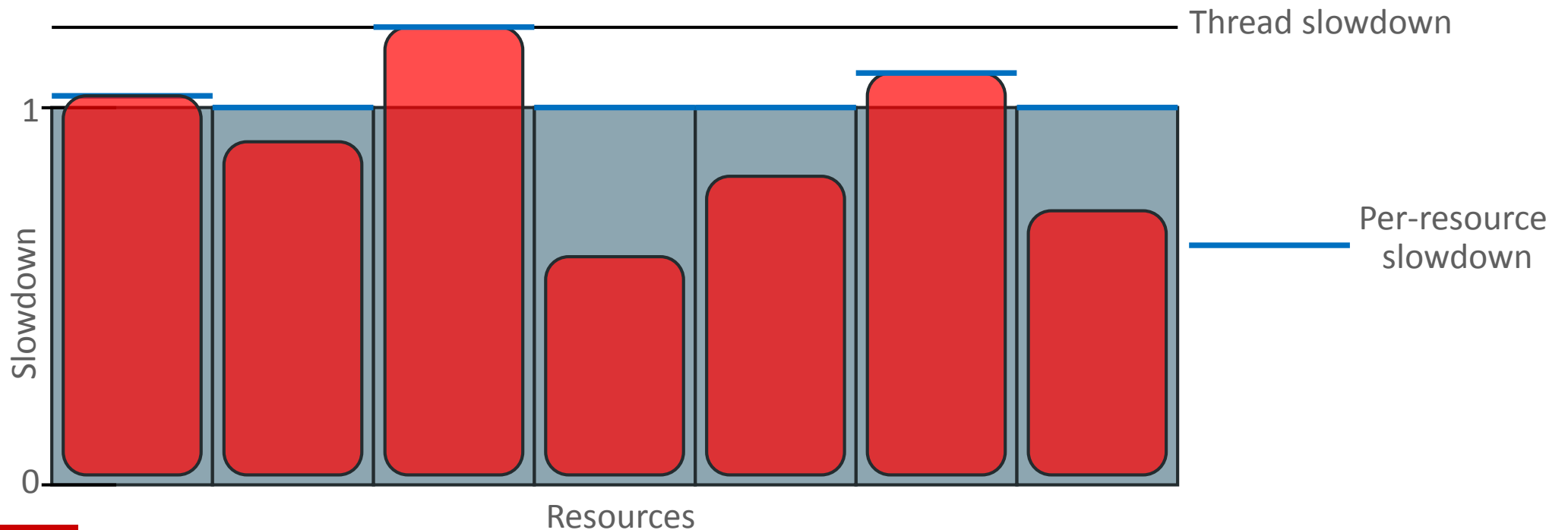
$$Slowdown = \frac{1.1}{1} = 1.1$$

ORACLE®

# Estimating thread slowdown

- Each thread has a slowdown calculated for all the resources it uses
- Each thread's slowdown is the maximum of the slowdowns it encounters

# Workload description: parallelism characteristics

- Characteristics of the interactions between the threads
- These all reflect synchronization either at the hardware level or within the application

| Parallelism | Communication slowdown |
|---|---|
| The percentage of the executed code that is parallel | Slowdown due to latency of thread communications between sockets, nodes, … |
| **Thread interlocking** | **Coincident resource demands** |
| Are threads independent? | Do threads have synchronized mode changes (e.g., intensive reads followed by intensive CPU) |

ORACLE®

# Workload description: parallelism characteristics

- Characteristics of the interactions between the threads
- These all reflect synchronization either at the hardware level or within the application
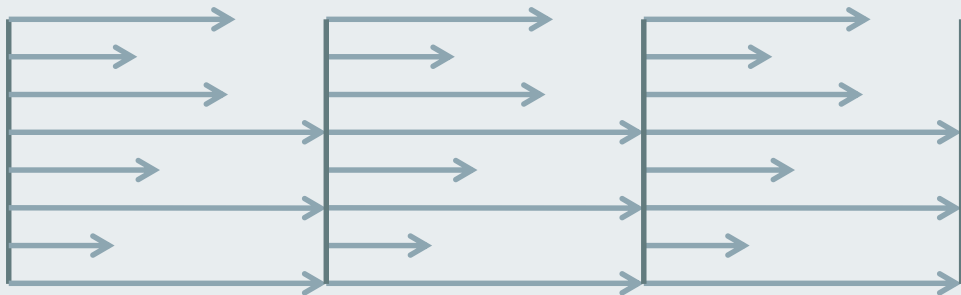
| | |
|---|---|
| **Parallelism**<br>The percentage of the executed code that is parallel | **Communication slowdown**<br>Slowdown due to latency of thread communications between sockets, nodes, … |
| **Thread interlocking**<br>Are threads independent? | **Coincident resource demands**<br>Do threads have synchronized mode changes (e.g., intensive reads followed by intensive CPU) |

ORACLE®

# Thread interlocking

- Describes how to combine thread slowdowns:

| Interlocked threads | Independent threads |
|---|---|
| Each thread does equal work | Each thread works for the same time |
|  |  |
| Return the lowest rate of work | Return the average rate of work |

# Thread interlocking



ep  swim  sp  sort join  is  applu  apsi  bwaves  pro  prh  md

ft  lu  mg  bt  wupwise  cg  npo  fma3d  prho  art  pagerank

Interlocked

Independent

# Thread interlocking



Run 2: Normal execution



Run 4: All threads slowed –
this provides the slowdown per thread

# Thread interlocking
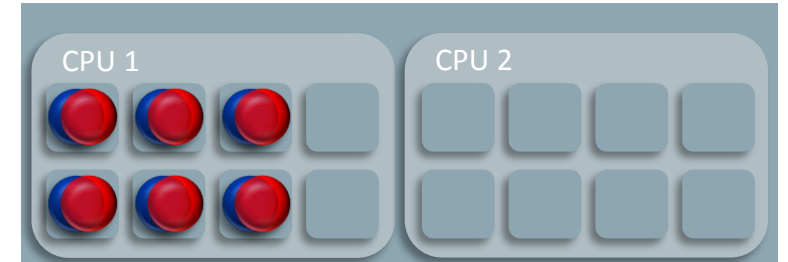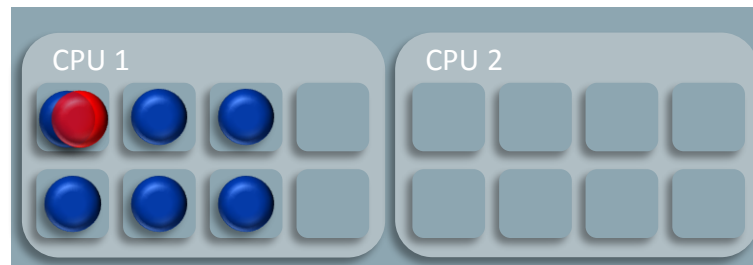


Run 2: Normal execution



Run 4: All threads slowed –
this provides the slowdown per thread



Run 5: One thread slowed –
this run's time is interpolated between the extremes

# Performance model

$$\text{Predicted speedup} = \frac{\text{Estimated speedup with thread count (Amdahl's law)}}{\text{Estimated slowdown due to resource contention}}$$

ORACLE®

# Thread utilization

**Workloads perform a roughly constant amount of work.**



Threads allocated

Resource demands for a unit of work with 1 thread

Time

Threads allocated

Ideal work rate with 2 threads.

Twice the resources for half the time

Time

Threads allocated

Realistic work rate with 2 threads.

Same total work, but spread more thinly over time.

Time

ORACLE®

# Contention a...

**Machine description generator**

Based on thread placement, predict demands on each resource (e.g., b/w on links in memory system).

Calculate per-thread slowdown based on their bottlenecks.

Add penalties for communication & synchronization.

Update estimates of thread utilization.

Proposed thread placement

**Performance prediction**

Speedup / sequential = 2.2

# Overview

1 ▸ What is the problem?

2 ▸ How does Pandia work?

**3 ▸ How well does Pandia work?**

4 ▸ Conclusions

# Evaluation

- Evaluated on a range of systems:
  - 2 socket Haswell, Ivy-bridge, Sandy-bridge
  - 4 socket Westmere

- Developed using 4 benchmarks (BT, CG, IS, MD), 18 additional benchmarks used in evaluation

- Profiles contain 9 parameters, but each test generates 1000s of data points

- Model features are tied to observable hardware and program features, not to features of the dataset

- Test the portability of workload descriptions between machines

ORACLE®

# Predicted v measured performance

# Predicted v measured performance

# Predicted v measured performance (Pagerank)

# Predicted v measured performance (Pagerank)

# Predicted v measured performance (Pagerank)



Increasing thread count

Measured
Predicted

Normalized speedup

Experiment placement number

ORACLE®

# Predicted v measured performance (Database hash join)



Chart: Normalized speedup (y-axis, 0.0 to 1.0) vs Experiment placement number (x-axis, 0 to 3000+). Arrow annotation: "Increasing thread count". Legend: Measured, Predicted.

# Average error

# Best vs predicted best placement

- Larger thread counts are more accurate

| Machine | Mean | Median |
|---|---|---|
| 2 socket Sandy-bridge | 0.77% | 0.00% |
| 2 socket Ivy-bridge | 0.29% | 0.00% |
| 2 socket Haswell | 2.78% | 1.05% |

ORACLE®

# Measured-best vs predicted-best placement (BT)

Measured-best



Predicted-best                                               Performance loss 0.36%

# Overview

1 ▶ What is the problem?

2 ▶ How does Pandia work?

3 ▶ How well does Pandia work?

4 ▶ Conclusions

# Conclusions

What we expected

| Modeling synchronization | Cache effects | NUMA effects |

What we found

ORACLE®

# Conclusions

Workloads can involve complex mixes of barriers, locks, atomics, etc.

**What we expected**

| Modeling synchronization | Cache effects | NUMA effects |

Averaging the effects with a simple model based on Amdahl's law was sufficient.

**What we found**

# Conclusions

|  |  |  |
|---|---|---|
| Workloads can involve complex mixes of barriers, locks, atomics, etc. | Source of complexity, and significant earlier exploration. | |
| **Modeling synchronization** | **Cache effects** | **NUMA effects** |
| Averaging the effects with a simple model based on Amdahl's law was sufficient. | Simple bandwidth based model suffices. More complex h/w mitigates performance cliffs. | |

*What we expected*

*What we found*

ORACLE®

# Conclusions

| | | |
|---|---|---|
| Workloads can involve complex mixes of barriers, locks, atomics, etc. | Source of complexity, and significant earlier exploration. | How to distinguish impact of local vs remote memory accesses? |
| **Modeling synchronization** | **Cache effects** | **NUMA effects** |
| Averaging the effects with a simple model based on Amdahl's law was sufficient. | Simple bandwidth based model suffices. More complex h/w mitigates performance cliffs. | Uniformity across the workload helps again: measure performance in aggregate. |

*What we expected*

*What we found*

# Conclusions

- Modern hardware avoids many of the pathological performance cases

- Simple models can be good enough to make meaningful decisions

- Predictions include resource predictions

- Best placements not always found by exploring scatter and pack placements

- State exploration will only get more complex when considering multiple workloads, so technique like Pandia are needed

For information about Pandia or roles in Oracle Labs please get in touch – daniel.goodman@oracle.com

**ORACLE®**

# Integrated Cloud

## Applications & Platform Services

**ORACLE**®